

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041

Deep Cauchy Hashing for Hamming Space Retrieval

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
Anonymous CVPR submission061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
Paper ID 1381

Abstract

Due to its computation efficiency and retrieval quality, hashing has been widely applied to approximate nearest neighbor search for large-scale image retrieval, while deep hashing further improves the retrieval quality by end-to-end representation learning and hash coding. With compact hash codes, Hamming space retrieval enables the most efficient constant-time search that returns data points within a given Hamming radius to each query, by hash table lookups instead of linear scan. However, subject to the weak capability of concentrating relevant images to be within a small Hamming ball due to mis-specified loss functions, existing deep hashing methods may underperform for Hamming space retrieval. This work presents Deep Cauchy Hashing (DCH), a novel deep hashing model that generates compact and concentrated binary hash codes to enable efficient and effective Hamming space retrieval. The main idea is to design a pairwise cross-entropy loss based on Cauchy distribution, which penalizes significantly on similar image pairs with Hamming distance larger than the given Hamming radius threshold. Comprehensive experiments demonstrate that DCH can generate highly concentrated hash codes and yield state-of-the-art Hamming space retrieval performance on three datasets, NUS-WIDE, CIFAR-10, and MS-COCO.

1. Introduction

In the big data era, large-scale and high-dimensional media data has been pervasive in search engines and social networks. To guarantee retrieval quality and computation efficiency, approximate nearest neighbor (ANN) search has attracted increasing attention. Parallel to the traditional indexing methods [16] for candidates pruning, another advantageous solution is hashing methods [29] for data compression, which transform high-dimensional media data into compact binary codes and generate similar binary codes for similar data items. This paper will focus on the learning to hash methods [29] that build data-dependent hash encoding schemes for efficient image retrieval, which have

shown better performance than the data-independent hashing methods, e.g. Locality-Sensitive Hashing (LSH) [8].

Many learning to hash methods have been proposed to enable efficient ANN search by Hamming ranking of compact hash codes, including both supervised and unsupervised methods [14, 10, 24, 7, 20, 28, 23, 22, 33]. Recently, deep learning to hash methods [31, 15, 26, 6, 34, 17, 19, 2] have shown that deep neural networks can enable end-to-end representation learning and hash coding with nonlinear hash functions. These deep learning to hash methods have shown state-of-the-art search performance. In particular, they prove it crucial to jointly learn similarity-preserving representations and control quantization error of converting continuous representations to binary codes [34, 17, 19, 2].

However, most of the existing methods are tailored to data compression instead of candidates pruning, i.e. they are designed to maximize retrieval performance based on linear scan of hash codes. As linear scan is still costly even using hash codes, we are somewhat deviating from our original intention with hashing, that is, to maximize speedup under acceptable accuracy. With the blossom of powerful hashing methods for linear scan, we should now turn to the Hamming space retrieval [7], which enables the most efficient constant-time search. In Hamming space retrieval, we return data points within a given Hamming radius to each query, by hash table lookups instead of linear scan. Unfortunately, existing hashing methods generally lack the capability of concentrating relevant images to be within a small Hamming ball due to the mis-specified loss functions, thus they may underperform for Hamming space retrieval.

This work presents Deep Cauchy Hashing (DCH), a novel deep hashing model that generates concentrated and compact hash codes to enable efficient and effective Hamming space retrieval. We propose a pairwise cross-entropy loss based on Cauchy distribution, which penalizes significantly on similar image pairs with Hamming distance larger than the given Hamming radius threshold. We further propose a quantization loss based on the Cauchy distribution, which enables learning nearly lossless hash codes. Both loss functions can be derived in a Bayesian learning framework and are well-specified to the Hamming space retrieval.

108 The proposed DCH model can be trained end-to-end by
 109 back-propagation. Extensive experiments demonstrate that
 110 DCH can generate highly concentrated and compact hash
 111 codes and yield state-of-the-art image retrieval performance
 112 on three datasets, NUS-WIDE, CIFAR-10, and MS-COCO.
 113

2. Related Work

Hashing Methods. Existing hashing methods [14, 10, 24, 7, 20, 28, 23, 9, 32, 33, 22] consist of unsupervised and supervised hashing. Please refer to [29] for a great survey.

Unsupervised hashing methods learn hash functions that encode data to binary codes by training from unlabeled data. Typical methods include reconstruction error minimization [25, 10, 12] and graph embedding [30, 21]. Supervised hashing further explores supervised information (e.g. pairwise similarity or relevance feedback) to generate discriminative and compact hash codes [14, 24, 20, 26]. Supervised Hashing with Kernels (KSH) [20] and Supervised Discrete Hashing (SDH) [26] generate nonlinear or discrete binary hash codes by minimizing (maximizing) the Hamming distances across similar (dissimilar) pairs of data points.

Recently, deep learning to hash methods [31, 15, 26, 6, 34, 17, 19, 2] yield breakthrough results on image retrieval datasets by blending the power of deep learning [13, 11]. In particular, DHN [34] is the first end-to-end framework that jointly preserves pairwise similarity and controls the quantization error. HashNet [2] improves DHN by balancing the positive and negative pairs in training data, and by continuation technique for lower quantization error, which obtains state-of-the-art performance on several benchmark datasets.

However, previous deep hashing methods perform unsatisfactorily for Hamming space retrieval [7], with early pruning to discard irrelevant data points out of Hamming Radius 2. The reason for inefficient Hamming space retrieval is that their loss functions penalize little when two similar data points have large Hamming distance. Thus they cannot concentrate relevant data points to be within Hamming radius 2. We propose a novel cross-entropy loss for similarity-preserving learning and a novel quantization loss for controlling hashing quality, both based on the Cauchy distribution. To our knowledge, this work is the first endeavor towards deep hashing for Hamming space retrieval.

3. Deep Cauchy Hashing

In similarity retrieval systems, we are given a training set of N points $\{\mathbf{x}_i\}_{i=1}^N$, each represented by a D -dimensional feature vector $\mathbf{x}_i \in \mathbb{R}^D$. Some pairs of points \mathbf{x}_i and \mathbf{x}_j are provided with similarity labels s_{ij} , where $s_{ij} = 1$ if \mathbf{x}_i and \mathbf{x}_j are similar while $s_{ij} = 0$ if \mathbf{x}_i and \mathbf{x}_j are dissimilar. Deep hashing learns a nonlinear hash function $f : \mathbf{x} \mapsto \mathbf{h} \in \{-1, 1\}^K$ from input space \mathbb{R}^D to Hamming space $\{-1, 1\}^K$ using deep neural networks, which encodes

each point \mathbf{x} into compact K -bit hash code $\mathbf{h} = f(\mathbf{x})$ such that the similarity information conveyed in given pairs \mathcal{S} can be preserved in the compact hash codes. In supervised hashing, the similarity information $\mathcal{S} = \{s_{ij}\}$ can be collected from the semantic labels of data points or relevance feedback from click-through data in online search engines.

Definition 1 (Hamming Space Retrieval). *For binary codes of K bits, the number of distinct hash buckets to examine is $N(K, r) = \sum_{k=0}^r \binom{K}{k}$, where r is the Hamming radius. $N(K, r)$ grows rapidly with r and when $r \leq 2$, it only requires $O(1)$ time for each query to find all r -neighbors. Hamming space retrieval refers to the retrieval scenario that directly returns data points within Hamming radius r to each query, by hash table lookups instead of linear scan.*

This paper presents a new Deep Cauchy Hashing (**DCH**) to enable efficient Hamming space retrieval in an end-to-end framework, as shown in Figure 1. The proposed deep architecture accepts pairwise input images $\{(\mathbf{x}_i, \mathbf{x}_j, s_{ij})\}$ and processes them through an end-to-end pipeline of deep representation learning and binary hash coding: (1) a convolutional network (CNN) for learning deep representation of each image \mathbf{x}_i , (2) a fully-connected hash layer (fch) for transforming the deep representation into K -bit hash code $\mathbf{h}_i \in \{1, -1\}^K$, (3) a novel Cauchy cross-entropy loss for similarity-preserving learning in Hamming space, and (4) a novel Cauchy quantization loss for controlling both the binarization error and the hashing quality in Hamming space.

3.1. Deep Architecture

The architecture for Deep Cauchy Hashing is shown in Figure 1. We extend from AlexNet [13], a deep convolutional neural network (CNN) with five convolutional layers $conv1$ – $conv5$ and three fully connected layers $fc6$ – $fc8$. We replace the classifier layer $fc8$ with a new hash layer fch of K hidden units, which transforms the representation of the $fc7$ layer into K -dimensional continuous code $\mathbf{z}_i \in \mathbb{R}^K$ for each image \mathbf{x}_i . We obtain hash code \mathbf{h}_i through the sign thresholding $\mathbf{h}_i = \text{sgn}(\mathbf{z}_i)$. However, since it is hard to optimize the sign function due to ill-posed gradient, we adopt the hyperbolic tangent (\tanh) function to squash the continuous code \mathbf{z}_i to be within $[-1, 1]$, which reduces the gap between the continuous code \mathbf{z}_i and the binary hash code \mathbf{h}_i . To further guarantee the quality of hash codes for efficient Hamming space retrieval, we preserve the similarity between the training pairs $\{(\mathbf{x}_i, \mathbf{x}_j, s_{ij}) : s_{ij} \in \mathcal{S}\}$ and control the quantization error, both performed in the Hamming space. Towards this goal, this paper proposes two novel loss functions based on the long-tailed Cauchy distribution: a pairwise Cauchy cross-entropy loss and a pointwise Cauchy quantization loss, both derived in the Maximum a Posteriori (MAP) estimation framework.

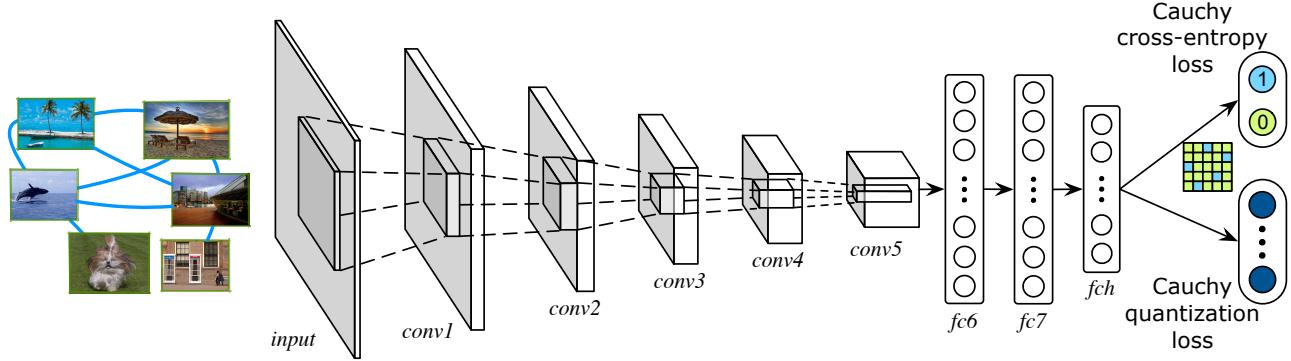


Figure 1. The architecture of the proposed Deep Cauchy Hashing (DCH), which is comprised of four key components: (1) a convolutional network (CNN) for learning deep representation of each image x_i , (2) a fully-connected hash layer (fch) for transforming the deep representation into K -bit hash code $\mathbf{h}_i \in \{1, -1\}^K$, (3) a novel Cauchy cross-entropy loss for similarity-preserving learning in the Hamming space, and (4) a novel Cauchy quantization loss for controlling both the binarization error and the hash code quality. *Best viewed in color.*

3.2. Bayesian Learning Framework

In this paper, we propose a Bayesian learning framework to perform deep hashing from similarity data by jointly preserving similarity of pairwise images and controlling the quantization error. Given training images with pairwise similarity labels as $\{(x_i, x_j, s_{ij}) : s_{ij} \in \mathcal{S}\}$, the logarithm Maximum a Posteriori (MAP) estimation of the hash codes $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N]$ for N training images can be defined as

$$\begin{aligned} \log P(\mathbf{H}|\mathcal{S}) &\propto \log P(\mathcal{S}|\mathbf{H}) P(\mathbf{H}) \\ &= \sum_{s_{ij} \in \mathcal{S}} w_{ij} \log P(s_{ij}|\mathbf{h}_i, \mathbf{h}_j) + \sum_{i=1}^N \log P(\mathbf{h}_i) \end{aligned} \quad (1)$$

where $P(\mathcal{S}|\mathbf{H}) = \prod_{s_{ij} \in \mathcal{S}} [P(s_{ij}|\mathbf{h}_i, \mathbf{h}_j)]^{w_{ij}}$ is the weighted likelihood function [4], and w_{ij} is the weight for each training pair (x_i, x_j, s_{ij}) , which tackles the data imbalance problem by weighting training pairs according to the importance of misclassifying that pair. Since similarity label can only be $s_{ij} = 1$ or $s_{ij} = 0$, to account for the data imbalance between similar and dissimilar pairs, we propose

$$w_{ij} = \begin{cases} |\mathcal{S}| / |\mathcal{S}_1|, & s_{ij} = 1 \\ |\mathcal{S}| / |\mathcal{S}_0|, & s_{ij} = 0 \end{cases} \quad (2)$$

where $\mathcal{S}_1 = \{s_{ij} \in \mathcal{S} : s_{ij} = 1\}$ is the set of similar pairs and $\mathcal{S}_0 = \{s_{ij} \in \mathcal{S} : s_{ij} = 0\}$ is the set of dissimilar pairs. For each pair, $P(s_{ij}|\mathbf{h}_i, \mathbf{h}_j)$ is the conditional probability of similarity label s_{ij} given a pair of hash codes \mathbf{h}_i and \mathbf{h}_j , which can be naturally defined by the Bernoulli distribution,

$$\begin{aligned} P(s_{ij}|\mathbf{h}_i, \mathbf{h}_j) &= \begin{cases} \sigma(d(\mathbf{h}_i, \mathbf{h}_j)), & s_{ij} = 1 \\ 1 - \sigma(d(\mathbf{h}_i, \mathbf{h}_j)), & s_{ij} = 0 \end{cases} \\ &= \sigma(d(\mathbf{h}_i, \mathbf{h}_j))^{s_{ij}} (1 - \sigma(d(\mathbf{h}_i, \mathbf{h}_j)))^{1-s_{ij}} \end{aligned} \quad (3)$$

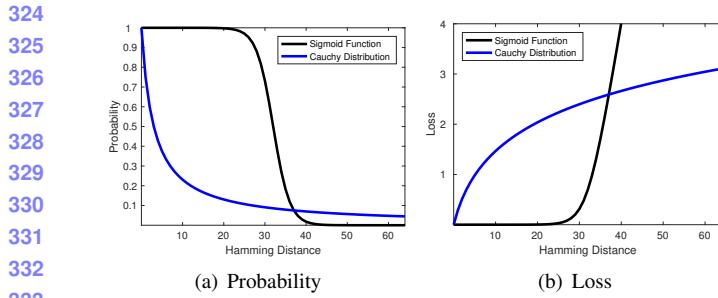
where $d(\mathbf{h}_i, \mathbf{h}_j)$ denotes the Hamming distance between hash codes \mathbf{h}_i and \mathbf{h}_j , and σ is a well-defined probability function to be elaborated in the next subsection.

Similar to binary-class logistic regression for pointwise data, we see in Equation (3) that the smaller the Hamming distance $d(\mathbf{h}_i, \mathbf{h}_j)$ is, the larger the conditional probability $P(1|\mathbf{h}_i, \mathbf{h}_j)$ will be, implying that the image pair x_i and x_j should be classified as similar; otherwise, the larger the conditional probability $P(0|\mathbf{h}_i, \mathbf{h}_j)$ will be, implying that the image pair should be classified as dissimilar. Hence, Equation (3) is a reasonable extension of the binary-class logistic regression to the pairwise classification scenario, which is a natural solution to the binary similarity labels $s_{ij} \in \{0, 1\}$.

3.3. Cauchy Hash Learning

With the Bayesian learning framework, any probability function σ and distance function d can be used to instantiate a specific hashing model. Previous state-of-the-art deep hashing methods, such as DHN [34] and HashNet [2], usually adopt sigmoid function $\sigma(x) = 1/(1 + e^{-x})$ as the probability function. However, we discover a key misspecification problem of the sigmoid function as illustrated in Figure 2. We can observe that the probability of sigmoid function stays high when the Hamming distance between hash codes is much larger than 2 and only starts to decrease obviously when the Hamming distance becomes close to $K/2$. This implies that previous deep hashing methods cannot pull the Hamming distance between the hash codes of similar data points to be smaller than 2, because the probabilities for different Hamming distances smaller than $K/2$ are not discriminative enough. This is a severe disadvantage of the existing hashing methods, which makes efficient Hamming space retrieval impossible. Note that in Hamming space retrieval, we can only return the data points within the Hamming ball of radius 2 for each query.

To tackle the above misspecification problem of sigmoid



(a) Probability (b) Loss

Figure 2. The values of Probability (a) and Loss (b) with respect to Hamming Distance between the hash codes of similar data points ($s_{ij} = 1$). The Probability (Loss) based on sigmoid function is very large (small) even for Hamming distance much larger than 2, which is ill-specified for Hamming ball retrieval. As a desired property, the Probability based on Cauchy distribution is large only for smaller Hamming distance; the Loss based on Cauchy distribution significantly penalizes for Hamming distance larger than 2.

function, we propose a novel probability function based on the Cauchy distribution with many nice properties as

$$\sigma(d(\mathbf{h}_i, \mathbf{h}_j)) = \frac{\gamma}{\gamma + d(\mathbf{h}_i, \mathbf{h}_j)}, \quad (4)$$

where γ is the scale parameter of the Cauchy distribution, and the normalization constant $\frac{1}{\pi\sqrt{\gamma}}$ is omitted for clarity, since it will not influence the final model. In Figure 2, we can observe that the probability of the proposed Cauchy distribution decreases very fast when the Hamming distance is small, resulting in that the similar points will be pulled to be within small Hamming radius. The decaying speed of the probability will be even faster by using a smaller γ , which imposes more force to concentrate similar points to be within small Hamming radius. Hence, the scale parameter γ is very important to control the tradeoff levels between precision and recall. By simply varying γ , we can support diverse Hamming space retrieval scenarios with different Hamming radiiuses to give different pruning rates.

Since discrete optimization of Equation (1) with binary constraints $\mathbf{h}_i \in \{-1, 1\}^K$ is very challenging, continuous relaxation is applied to the binary constraints for ease of optimization, as adopted by most previous hashing methods [29, 34]. To control the quantization error $\|\mathbf{h}_i - \text{sgn}(\mathbf{h}_i)\|$ caused by continuous relaxation and to learn high-quality hash codes, we propose a novel prior for each hash code \mathbf{h}_i based on a symmetric variant of the Cauchy distribution as

$$P(\mathbf{h}_i) = \frac{\gamma}{\gamma + d(|\mathbf{h}_i|, \mathbf{1})}, \quad (5)$$

where γ is the scale parameter of the symmetric Cauchy distribution, and $\mathbf{1} \in \mathbb{R}^K$ is the vector of ones.

By using the continuous relaxation, we need to replace the Hamming distance with its best approximation on continuous codes. For a pair of binary hash codes \mathbf{h}_i and \mathbf{h}_j ,

there exists a nice relationship between their Hamming distance $d(\mathbf{h}_i, \mathbf{h}_j)$ and the normalized Euclidean distance as

$$\begin{aligned} d(\mathbf{h}_i, \mathbf{h}_j) &= \frac{K}{4} \left\| \frac{\mathbf{h}_i}{\|\mathbf{h}_i\|} - \frac{\mathbf{h}_j}{\|\mathbf{h}_j\|} \right\|_2^2 \\ &= \frac{K}{2} (1 - \cos(\mathbf{h}_i, \mathbf{h}_j)). \end{aligned} \quad (6)$$

Hence this paper adopts $d(\mathbf{h}_i, \mathbf{h}_j) = \frac{K}{2} (1 - \cos(\mathbf{h}_i, \mathbf{h}_j))$. After using the continuous relaxation to ease the optimization, the normalized Euclidean distance compares each pair of continuous codes on a unit sphere, while the Hamming distance compares each pair of hash codes on a unit hypercube. More intuitively, the unit sphere is always a circumscribed sphere of the hypercube. As a desirable result, by mitigating the diversity of code lengths, the normalized Euclidean distance on continuous codes is always the upper bound of the Hamming distance on hash codes.

By taking Equations (3) and (5) into the MAP estimation framework in Equation (1), we obtain the optimization problem of the proposed Deep Cauchy Hashing (DCH) as

$$\min_{\Theta} L + \lambda Q, \quad (7)$$

where λ is a hyper-parameter to trade-off the Cauchy cross-entropy loss L and the Cauchy quantization loss Q , and Θ denotes the set of network parameters to be optimized. Specifically, the Cauchy cross-entropy loss L is derived as

$$L = \sum_{s_{ij} \in \mathcal{S}} w_{ij} \left(s_{ij} \log \frac{d(\mathbf{h}_i, \mathbf{h}_j)}{\gamma} + \log \left(1 + \frac{\gamma}{d(\mathbf{h}_i, \mathbf{h}_j)} \right) \right), \quad (8)$$

and similarly, the Cauchy quantization loss is derived as

$$Q = \sum_{i=1}^N \log \left(1 + \frac{d(|\mathbf{h}_i|, \mathbf{1})}{\gamma} \right), \quad (9)$$

where $d(\cdot, \cdot)$ is either the Hamming distance between the hash codes or the normalized Euclidean distance between the continuous codes. Since the quantization error will be controlled by the proposed Cauchy quantization loss in the joint optimization problem (7), for ease of optimization, we can use continuous relaxation for the hash codes \mathbf{h}_i during training (only for training, not for testing).

Based on the MAP estimation in Equation (7), we can enable statistically optimal learning of compact hash codes by jointly preserving the pairwise similarity and controlling the quantization error. Finally, we can obtain K -bit binary codes by simple sign thresholding $\mathbf{h} \leftarrow \text{sgn}(\mathbf{h})$, where $\text{sgn}(\mathbf{h})$ is the sign function on vectors that for $i = 1, \dots, K$, $\text{sgn}(h_i) = 1$ if $h_i > 0$, otherwise $\text{sgn}(h_i) = -1$. It is worth noting that, since we have minimized the quantization error in (7) during training, this final binarization step will incur very small loss of retrieval quality as validated empirically.

432

4. Experiments

433

We conduct extensive experiments to evaluate the efficacy of the proposed DCH approach with several state-of-the-art shallow and deep hashing methods on three benchmark datasets. Codes and configurations will be released.

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

4.1. Setup

The evaluation is conducted on three widely-used benchmark datasets: NUS-WIDE, MS-COCO and CIFAR-10.

NUS-WIDE¹ [3] is a benchmark dataset that contains 269,648 images from Flickr.com. Each image is manually annotated by some of the 81 ground truth concepts (categories) for evaluating retrieval models. We follow similar experimental protocols in [34, 2], and randomly sample 5,000 images as the query points, with the remaining images used as the database and randomly sample 10,000 images from the database as the training points.

MS-COCO² [18] is a popular dataset for image recognition, segmentation and captioning. The current release contains 82,783 training images and 40,504 validation images, where each image is labeled by some of the 80 semantic concepts. We randomly sample 5,000 images as query points, with the rest used as the database, and randomly sample 10,000 images from the database for training.

CIFAR-10³ is a standard dataset with 60,000 images in 10 classes. We follow protocol in [34, 1] to randomly select 100 images per class as query set, 500 images per class as training set, and the rest images are used as the database.

Following standard protocol as in [31, 15, 34, 2], the similarity information for hash learning and for ground-truth evaluation is constructed from image labels: if two images i and j share at least one label, they are similar and $s_{ij} = 1$; otherwise, they are dissimilar and $s_{ij} = 0$. Note that, although we use the image labels to construct the similarity information, the proposed approach DCH can learn hash codes when only the similarity information is available. By constructing the training data in this way, the ratio between the number of dissimilar pairs and the number of similar pairs is roughly 10, 5, and 1 for CIFAR-10, NUS-WIDE, and MS-COCO, respectively. These datasets exhibit the data imbalance phenomenon and can be used to evaluate different hashing methods under data imbalance scenario.

We compare the retrieval performance of **DCH** with eight classical or state-of-the-art hashing methods: supervised shallow methods **ITQ-CCA** [10], **BRE** [14], **KSH** [20], and **SDH** [26], and supervised deep methods **CNNH** [31], **DNNH** [15], **DHN** [34], and **HashNet** [2].

We follow standard evaluation methods for Hamming space retrieval [7], which consists of two consecutive steps:

¹<http://lms.comp.nus.edu.sg/research/NUS-WIDE.htm>

²<http://mscoco.org>

³<http://www.cs.toronto.edu/kriz/cifar.html>

(1) Pruning, to return data points within Hamming radius 2 for each query using hash table lookups; (2) Scanning, to re-rank the returned data points in ascending order of their distances to each query using continuous codes. To evaluate the effectiveness of Hamming space retrieval, we report three standard evaluation metrics to measure the quality of the data points within Hamming radius 2: Mean Average Precision within Hamming Radius 2 (**MAP@H \leq 2**), Precision curves within Hamming Radius 2 (**P@H \leq 2**), and Recall curves within Hamming Radius 2 (**R@H \leq 2**).

Due to the potentially best efficiency, the search based on pruning followed by re-ranking is widely-deployed in large-scale online retrieval systems such as search engines. In the **MAP@H \leq 2** evaluation metric, we compute the Mean Average Precision for the re-ranked list of the data points pruned by Hamming radius 2, which can measure the search quality. More specifically, given a set of queries, we first compute the Average Precision (AP) of each query as

$$\text{AP}@T = \frac{\sum_{t=1}^T P(t) \delta(t)}{\sum_{t'=1}^T \delta(t')}, \quad (10)$$

where T is the number of top-returned data points within Hamming radius 2, $P(t)$ denotes the precision of top t retrieved results, and $\delta(t) = 1$ if the t -th retrieved result is a true neighbor of the query, otherwise $\delta(t) = 0$. Then **MAP@H \leq 2** is computed as the mean of average precisions for all queries. The larger the **MAP@H \leq 2**, the better the search quality for the data points within Hamming radius 2.

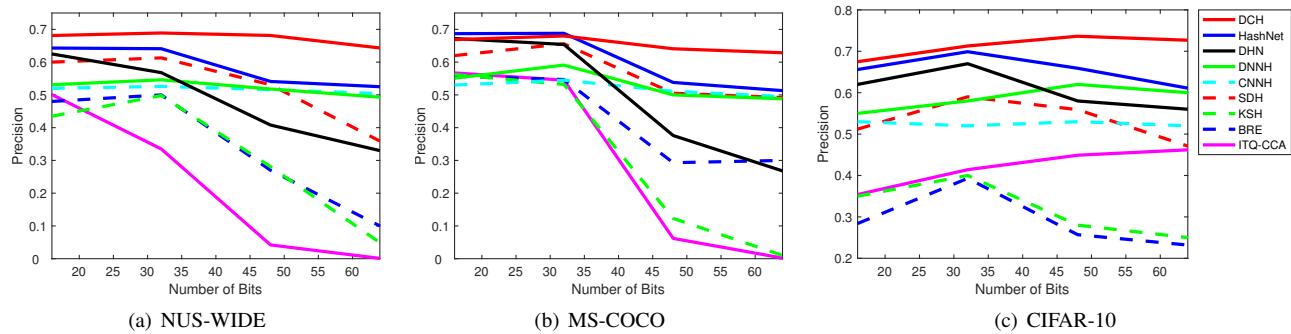
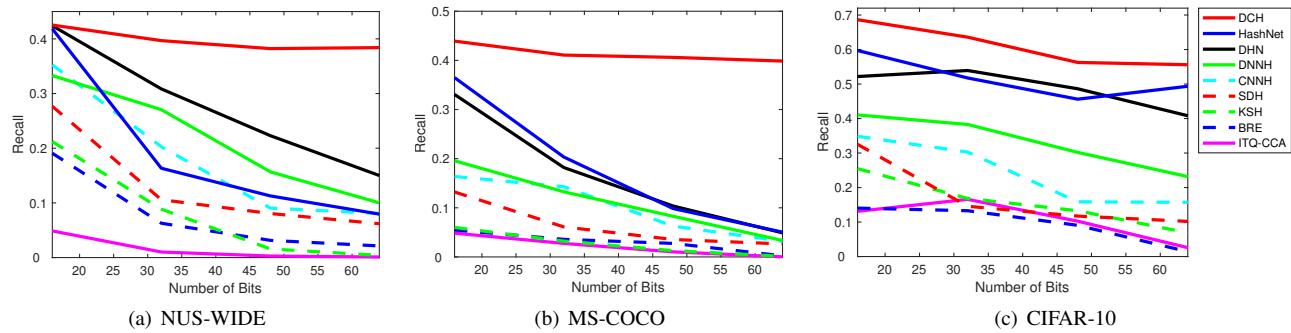
For shallow hashing methods, we use as image features the 4096-dimensional DeCAF₇ features [5]. For deep hashing methods, we use raw images as the input. We adopt the AlexNet architecture [13] for all deep hashing methods, and implement DCH based on the **TensorFlow** framework. We fine-tune convolutional layers *conv1*–*conv5* and fully-connected layers *fc6*–*fc7* copied from the AlexNet model pre-trained on ImageNet and train the last layer, all through back-propagation. As the last layer is trained from scratch, we set its learning rate to be 10 times that of the lower layers. We use mini-batch stochastic gradient descent (SGD) with 0.9 momentum and cross-validate the learning rate from 10^{-5} to 10^{-2} with a multiplicative step-size $10^{\frac{1}{2}}$. We fix the mini-batch size of images as 256 and the weight decay parameter as 0.0005. We select the model parameters of DCH, λ and γ , by cross-validation. We also select the parameters of each comparison method by cross-validation.

4.2. Results

The Mean Average Precision of Re-ranking within Hamming Radius 2 (**MAP@H \leq 2**) results of all comparison methods are listed in Table 1. Results show that DCH substantially outperforms all comparison methods, in that DCH can perform high-quality pruning within Hamming radius 2

540 Table 1. Mean Average Precision of Re-ranking within Hamming Radius 2 (MAP@H \leq 2) for Different Bits on Three Benchmark Datasets 594

541 Method	542 NUS-WIDE				543 MS-COCO				544 CIFAR-10			
	545 16 bits	546 32 bits	547 48 bits	548 64 bits	549 16 bits	550 32 bits	551 48 bits	552 64 bits	553 16 bits	554 32 bits	555 48 bits	556 64 bits
ITQ-CCA [10]	0.5706	0.4397	0.0825	0.0051	0.5949	0.5612	0.0585	0.0105	0.4258	0.4652	0.4774	0.4932
BRE [14]	0.5502	0.5422	0.4128	0.2202	0.5625	0.5498	0.4214	0.4014	0.4216	0.4519	0.4002	0.3438
KSH [20]	0.5185	0.5659	0.4102	0.0608	0.5797	0.5532	0.2338	0.0216	0.4368	0.4585	0.4012	0.3819
SDH [26]	0.6681	0.6824	0.5979	0.4679	0.6449	0.6766	0.5226	0.5108	0.5620	0.6428	0.6069	0.5012
CNNH [31]	0.5843	0.5989	0.5734	0.5729	0.5602	0.5685	0.5376	0.5058	0.5512	0.5468	0.5454	0.5364
DNNH [15]	0.6191	0.6216	0.5902	0.5626	0.5771	0.6023	0.5235	0.5013	0.5703	0.5985	0.6421	0.6118
DHN [34]	0.6901	0.7021	0.6685	0.5664	0.6749	0.6680	0.5151	0.4186	0.6929	0.6445	0.5835	0.5883
HashNet [2]	0.6944	0.7147	0.6736	0.6190	0.6851	0.6900	0.5589	0.5344	0.7476	0.7776	0.6399	0.6259
DCH	0.7401	0.7720	0.7685	0.7124	0.7010	0.7576	0.7251	0.7013	0.7901	0.7979	0.8071	0.7936

563 Figure 3. The Precision curves within Hamming Radius 2 (P@H \leq 2) of DCH and comparison methods on the three benchmark datasets. 617576 Figure 4. The Recall curves within Hamming Radius 2 (R@H \leq 2) of DCH and comparison methods on the three benchmark datasets. 630

578 in the first step, enabling efficient re-ranking in the second 579 step. Specifically, compared to SDH, the best shallow hashing 580 method with deep features as input, DCH achieves 581 absolute increases of **14.4%**, **13.3%** and **21.9%** in average 582 MAP@H \leq 2 with different code lengths on NUS-WIDE, 583 MS-COCO and CIFAR-10, respectively. DCH outperforms 584 HashNet, the state-of-the-art deep hashing method, by large 585 margins of **7.3%**, **10.4%** and **9.9%** in average MAP@H \leq 2 586 with different code lengths on three benchmark datasets. 587

The MAP@H \leq 2 results reveal some interesting insights. 588 **(1)** Shallow hashing methods cannot learn discriminative 589 deep features and compact hash codes through end-to-end 590 framework, which explains the fact that they are surpassed 591 by deep hashing methods. **(2)** Deep hashing methods DHN 592 and HashNet learn less lossy hash codes by preserving 593 the similarity information and controlling the quantization

error, which also significantly outperform deep methods CNNH and DNNH without reducing the quantization error. 632

The proposed DCH model improves substantially from 633 the state-of-the-art HashNet by two important perspectives: 634 **(1)** DCH preserves similarity relationships based on Cauchy 635 distribution, which can achieve better pruning performance 636 within Hamming radius 2; **(2)** DCH learns the novel Cauchy 637 cross-entropy loss and Cauchy quantization loss based on 638 normalized distance, which can better approximate the 639 Hamming distance to learn nearly lossless hash codes. 640

The performance of Precision within Hamming Radius 641 2 (P@H \leq 2) is very important for Hamming space retrieval, 642 since it only requires $O(1)$ time for each query and enables 643 really efficient pruning. As shown in Figure 3, DCH often 644 achieves the highest P@H \leq 2 results on all three benchmark 645 datasets with regard to different code lengths. This 646

648 Table 2. Mean Average Precision of Re-ranking within Hamming Radius 2 (MAP@H≤2) of DCH and Its Variants on Three Datasets 702
649

Method	NUS-WIDE				MS-COCO				CIFAR-10			
	16 bits	32 bits	48 bits	64 bits	16 bits	32 bits	48 bits	64 bits	16 bits	32 bits	48 bits	64 bits
DCH	0.7401	0.7720	0.7685	0.7124	0.7010	0.7576	0.7251	0.7013	0.7901	0.7979	0.8071	0.7936
DCH-Q	0.7086	0.7458	0.7432	0.7028	0.6785	0.7238	0.6982	0.6913	0.7645	0.7789	0.7858	0.7832
DCH-C	0.6997	0.7205	0.6874	0.6328	0.6767	0.6972	0.5801	0.5536	0.7513	0.7628	0.6819	0.6627
DCH-E	0.7178	0.7511	0.7302	0.6982	0.6826	0.7128	0.6803	0.6735	0.7598	0.7739	0.7495	0.7287

656 validates that DCH can learn compacter hash codes than all
657 comparison methods and can enable more efficient and ac-
658 curate Hamming space retrieval. With longer hash codes,
659 the Hamming space will become more sparse and fewer
660 data points will fall in the Hamming ball within radius 2.
661 This is why most previous hashing methods achieve worse
662 retrieval performance with longer code lengths. It is worth
663 noting that DCH achieves a relatively mild decrease or even
664 an increase in accuracy by longer code lengths, validating
665 that DCH can effectively concentrate hash codes of simi-
666 lar data points together to be within the Hamming radius 2,
667 which significantly benefit Hamming space retrieval.

668 The performance of Recall within Hamming Radius 2
669 ($R@H\leq 2$) is crucial for Hamming space retrieval, since it
670 is likely that all data points will be pruned out due to the
671 highly sparse Hamming space. As shown in Figure 4, DCH
672 achieves the highest $R@H\leq 2$ results on all three datasets
673 w.r.t different code lengths, which is very encouraging. This
674 validates that DCH can concentrate more relevant points to
675 be within the Hamming ball of radius 2 than all the com-
676 parison methods. As the Hamming space will become more
677 sparse when using longer hash codes, most hashing base-
678 lines incur serious performance drop on $R@H\leq 2$. Since the
679 relationships between data pairs on multi-label datasets are
680 more complex than that on single-label datasets, the per-
681 formance drop becomes more serious on multi-label datasets
682 such as NUS-WIDE and MS-COCO. By introducing the
683 novel Cauchy cross-entropy loss and Cauchy quantization
684 loss, even on multi-label datasets, the proposed DCH incurs
685 very small performance drop on $R@H\leq 2$ as the hash codes
686 become longer, showing that DCH can concentrate more
687 relevant points to be within Hamming radius 2 even using
688 longer code lengths. The ability to use longer codes gives a
689 flexible to tradeoff between accuracy and efficiency, a flexi-
690 bility that is often impossible for previous hashing methods.
691

692 4.3. Discussion

693 4.3.1 Ablation Study

695 We investigate three DCH variants: (1) **DCH-Q** is a DCH
696 variant without using the new Cauchy quantization loss (9),
697 in other words $\lambda=0$; (2) **DCH-C** is a DCH variant replac-
698 ing the Cauchy cross-entropy loss with the popular sigmoid
699 cross-entropy loss [34, 2]; (3) **DCH-E** is a DCH variant
700 replacing the normalized Euclidean distance (6) with the
701 Euclidean distance as $d(\mathbf{h}_i, \mathbf{h}_j) = \frac{1}{4} \|\mathbf{h}_i - \mathbf{h}_j\|_2^2$. The

710 MAP@H≤2 results w.r.t. different code lengths on all three
711 benchmark datasets are reported in Table 2.

712 **Cauchy Cross-Entropy Loss.** DCH outperforms DCH-
713 C by very large margins of 6.3%, 9.4% and 8.3% in average
714 MAP@H≤2 with different code lengths on NUS-WIDE,
715 MS-COCO and CIFAR-10, respectively. The Cauchy cross-
716 entropy loss (8) is based on the Cauchy distribution, which
717 can keep more relevant points to be within small Hamming
718 radius to enable effective Hamming space retrieval, whereas
719 the sigmoid cross-entropy loss cannot achieve this desired
720 property. Also, DCH outperforms DCH-E by large margins
721 of 2.4%, 3.4% and 4.4% in average MAP@H≤2 with dif-
722 ferent code lengths on three datasets. In real search engines,
723 normalized distance is widely used to mitigate the diversity
724 of vector lengths and improve the retrieval quality, which
725 has not been integrated with the cross-entropy loss [29].

726 **Cauchy Quantization Loss.** DCH outperforms DCH-
727 Q by 2.3%, 2.3%, and 1.9% in average MAP@H≤2 with
728 different code lengths on three benchmarks, respectively.
729 These results validate that the novel Cauchy quantization
730 loss (9) can enhance the pruning efficiency and improve the
731 pruning and re-ranking results within Hamming radius 2.
732

733 4.3.2 Sensitivity Study

735 In Hamming space retrieval, a key aspect is to control the
736 tradeoff between precision and recall as well as efficiency
737 w.r.t. different Hamming radiiuses. As aforementioned, the
738 time cost for Hamming pruning through hash table lookups
739 is $N(K, r) = \sum_{k=0}^r \binom{K}{k}$, where r is the Hamming radius.
740 Hence $N(K, r) \propto K^r$, and we can only tolerate smaller
741 Hamming radius, typically $r \leq 10$. In other words, we
742 cannot use large Hamming radius for higher recall. In this
743 paper, we introduce the Cauchy distribution parameter γ to
744 fully tradeoff precision and recall, with sensitivity study of
745 γ in terms of precision and recall shown in Figure 5.

746 Figure 5(a) shows the precision curves w.r.t. different
747 values of γ for practical Hamming radiiuses $r = [0, \dots, 4]$.
748 Figure 5(b) shows the precision curves w.r.t. different Ham-
749 ming radiiuses for typical values of $\gamma = [2, \dots, 500]$, which
750 is an alternative view of Figure 5(a). As can be seen, larger
751 (smaller) Hamming radius requires larger (smaller) value
752 of γ to guarantee the highest precision, which is consistent
753 with the theoretical analysis. Although larger Hamming ra-
754 dius leads to higher precision when larger γ is used, the
755 pruning cost will be exponentially enlarged as $O(K^r)$.

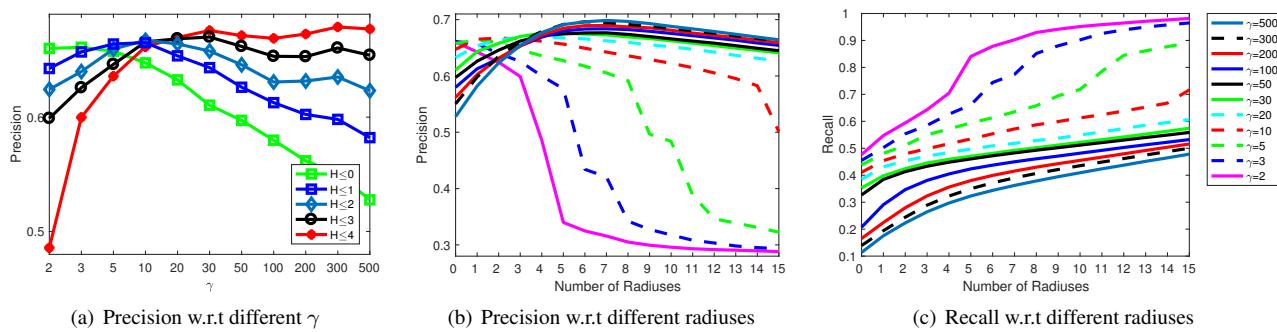


Figure 5. Sensitivity study for DCH using 64-bit hash codes on NUS-WIDE dataset: (a) Precision curves w.r.t. γ for different Hamming radiiuses, (b) Precision curves w.r.t. Hamming radiiuses for different γ , and (c) Recall curves w.r.t. Hamming radiiuses for different γ .

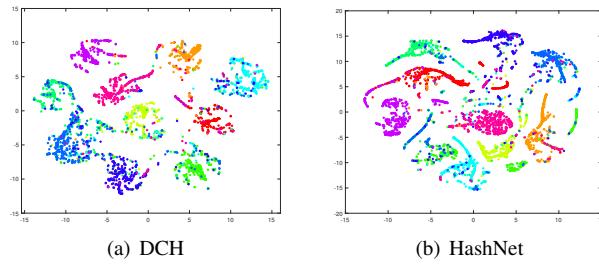


Figure 6. The t-SNE visualization of hash codes on CIFAR-10.

Figure 5(c) shows the recall curves w.r.t. different Hamming radiiuses for typical values of $\gamma = [2, \dots, 500]$. As can be seen, larger (smaller) Hamming radius leads to higher (lower) recall, which is consistent with the theoretical analysis. A very interesting observation is that, smaller (larger) γ leads to larger (smaller) recall. As analyzed before, γ controls the power of concentrating relevant data points within smaller Hamming balls, and smaller γ leads to larger concentration power. Again, although larger Hamming radius leads to higher recall when smaller γ is used, the pruning cost will be exponentially enlarged as $O(K^r)$.

Therefore, DCH provides with the powerful flexibility to tradeoff precision and recall as well the pruning efficiency. In practice, we usually decide the Hamming radius first based on the efficiency requirement, and typically $r = 2$. Then we tradeoff the precision and recall solely based on γ . As seen from Figure 5, $\gamma = 5$ is the best choice for Hamming radius $r = 2$. Besides the optimal value, we can vary $\gamma \in [2, 50]$ to achieve both satisfactory precision and recall.

4.3.3 Visualization Analysis

Visualization of Hash Codes by t-SNE. Figure 6 shows the t-SNE visualization [27] of the hash codes learned by DCH and the best deep hashing baseline HashNet [2] on CIFAR-10 dataset. We can observe that the hash codes generated by DCH show clear discriminative structures where the hash codes in different categories are well separated, while the hash codes generated by HashNet do not show

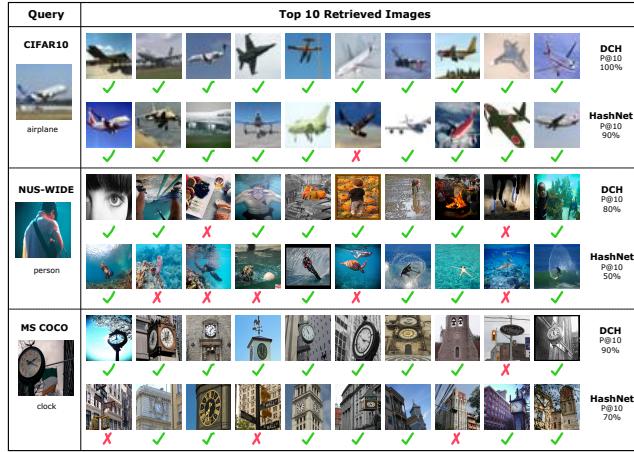


Figure 7. The top 10 images returned by DCH and HashNet.

such clear structures. This verifies that by introducing the Cauchy distribution for hashing, the hash codes generated through DCH are more discriminative than that generated by HashNet, enabling more effective image retrieval.

Illustration of Top 10 Returned Images. In Figure 7, as an intuitive illustration, we visualize the top 10 returned images of DCH and the best deep hashing baseline HashNet [2] for three query images on NUS-WIDE, MS-COCO and CIFAR-10, respectively. It shows that DCH can yield much more relevant and user-desired retrieval results.

5. Conclusion

This paper establishes efficient and effective Hamming space retrieval with constant-time search complexity. The proposed Deep Cauchy Hashing (DCH) approach generates compact and concentrated hash codes by jointly optimizing a novel Cauchy cross-entropy loss and a Cauchy quantization loss in a single Bayesian learning framework. The overall model can be trained end-to-end with well-specified loss functions. Extensive experiments show that DCH can yield state-of-the-art Hamming space retrieval performance on three datasets, NUS-WIDE, CIFAR-10, and MS-COCO.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

References

- [1] Y. Cao, M. Long, J. Wang, H. Zhu, and Q. Wen. Deep quantization network for efficient image retrieval. In *AAAI*. AAAI, 2016. 5
- [2] Z. Cao, M. Long, J. Wang, and P. S. Yu. Hashnet: Deep learning to hash by continuation. *arXiv preprint arXiv:1702.00758*, 2017. 1, 2, 3, 5, 6, 7, 8
- [3] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y.-T. Zheng. Nus-wide: A real-world web image database from national university of singapore. In *ICMR*. ACM, 2009. 5
- [4] J. P. Dmochowski, P. Sajda, and L. C. Parra. Maximum likelihood in cost-sensitive learning: Model specification, approximations, and upper bounds. *Journal of Machine Learning Research (JMLR)*, 11(Dec):3313–3332, 2010. 3
- [5] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*, 2014. 5
- [6] V. Erin Liang, J. Lu, G. Wang, P. Moulin, and J. Zhou. Deep hashing for compact binary codes learning. In *CVPR*, pages 2475–2483. IEEE, 2015. 1, 2
- [7] D. J. Fleet, A. Punjani, and M. Norouzi. Fast search in hamming space with multi-index hashing. In *CVPR*. IEEE, 2012. 1, 2, 5
- [8] A. Gionis, P. Indyk, R. Motwani, et al. Similarity search in high dimensions via hashing. In *VLDB*, volume 99, pages 518–529. ACM, 1999. 1
- [9] Y. Gong, S. Kumar, H. Rowley, S. Lazebnik, et al. Learning binary codes for high-dimensional data using bilinear projections. In *CVPR*, pages 484–491. IEEE, 2013. 2
- [10] Y. Gong and S. Lazebnik. Iterative quantization: A procrustean approach to learning binary codes. In *CVPR*, pages 817–824, 2011. 1, 2, 5, 6
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CVPR*, 2016. 2
- [12] H. Jegou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 33(1):117–128, Jan 2011. 2
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 2, 5
- [14] B. Kulis and T. Darrell. Learning to hash with binary reconstructive embeddings. In *NIPS*, pages 1042–1050, 2009. 1, 2, 5, 6
- [15] H. Lai, Y. Pan, Y. Liu, and S. Yan. Simultaneous feature learning and hash coding with deep neural networks. In *CVPR*. IEEE, 2015. 1, 2, 5, 6
- [16] M. S. Lew, N. Sebe, C. Djerafa, and R. Jain. Content-based multimedia information retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 2(1):1–19, Feb. 2006. 1
- [17] W.-J. Li, S. Wang, and W.-C. Kang. Feature learning based deep supervised hashing with pairwise labels. In *IJCAI*, 2016. 1, 2
- [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. 5
- [19] H. Liu, R. Wang, S. Shan, and X. Chen. Deep supervised hashing for fast image retrieval. In *CVPR*, pages 2064–2072, 2016. 1, 2
- [20] W. Liu, J. Wang, R. Ji, Y.-G. Jiang, and S.-F. Chang. Supervised hashing with kernels. In *CVPR*. IEEE, 2012. 1, 2, 5, 6
- [21] W. Liu, J. Wang, S. Kumar, and S.-F. Chang. Hashing with graphs. In *ICML*. ACM, 2011. 2
- [22] X. Liu, J. He, C. Deng, and B. Lang. Collaborative hashing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2139–2146, 2014. 1, 2
- [23] X. Liu, J. He, B. Lang, and S.-F. Chang. Hash bit selection: a unified solution for selection problems in hashing. In *CVPR*, pages 1570–1577. IEEE, 2013. 1, 2
- [24] M. Norouzi and D. M. Blei. Minimal loss hashing for compact binary codes. In *ICML*, pages 353–360. ACM, 2011. 1, 2
- [25] R. Salakhutdinov and G. E. Hinton. Learning a nonlinear embedding by preserving class neighbourhood structure. In *AISTATS*, pages 412–419, 2007. 2
- [26] F. Shen, C. Shen, W. Liu, and H. Tao Shen. Supervised discrete hashing. In *CVPR*. IEEE, June 2015. 1, 2, 5, 6
- [27] L. van der Maaten and G. Hinton. Visualizing high-dimensional data using t-sne. *Journal of Machine Learning Research*, 9: 2579–2605, Nov 2008. 8
- [28] J. Wang, S. Kumar, and S.-F. Chang. Semi-supervised hashing for large-scale search. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 34(12):2393–2406, 2012. 1, 2
- [29] J. Wang, H. T. Shen, J. Song, and J. Ji. Hashing for similarity search: A survey. Arxiv, 2014. 1, 2, 4, 7
- [30] Y. Weiss, A. Torralba, and R. Fergus. Spectral hashing. In *NIPS*, 2009. 2
- [31] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan. Supervised hashing for image retrieval via image representation learning. In *AAAI*, pages 2156–2162. AAAI, 2014. 1, 2, 5, 6
- [32] F. X. Yu, S. Kumar, Y. Gong, and S.-F. Chang. Circulant binary embedding. In *ICML*, pages 353–360. ACM, 2014. 2
- [33] P. Zhang, W. Zhang, W.-J. Li, and M. Guo. Supervised hashing with latent factor models. In *SIGIR*, pages 173–182. ACM, 2014. 1, 2
- [34] H. Zhu, M. Long, J. Wang, and Y. Cao. Deep hashing network for efficient similarity retrieval. In *AAAI*. AAAI, 2016. 1, 2, 3, 4, 5, 6, 7